

Реконструкция последовательности белка с учетом вариаций

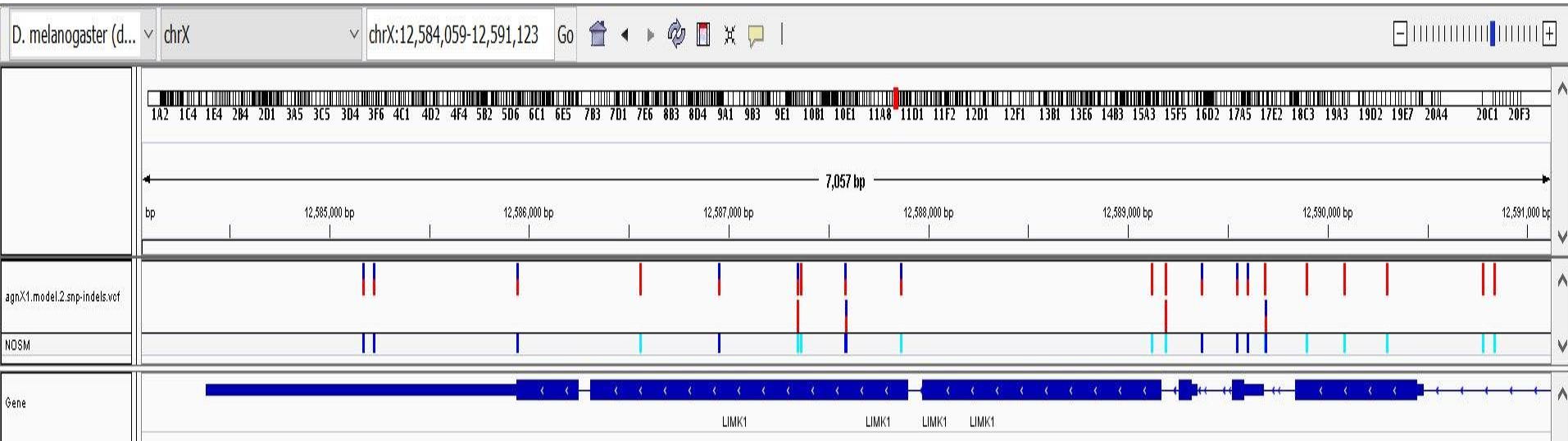
Научный руководитель: Захаров Геннадий,

Участники проекта: Нестеренко Максим, Цуринов Петр,
Кожевникова Ольга, Коляденко Илья

Научные консультанты: Зуева Мария, Альперович Михаил

Цель работы

Создать инструмент для восстановления аминокислотной последовательности белка, с учетом всех мутаций, локализованных в белок-кодирующих участках.

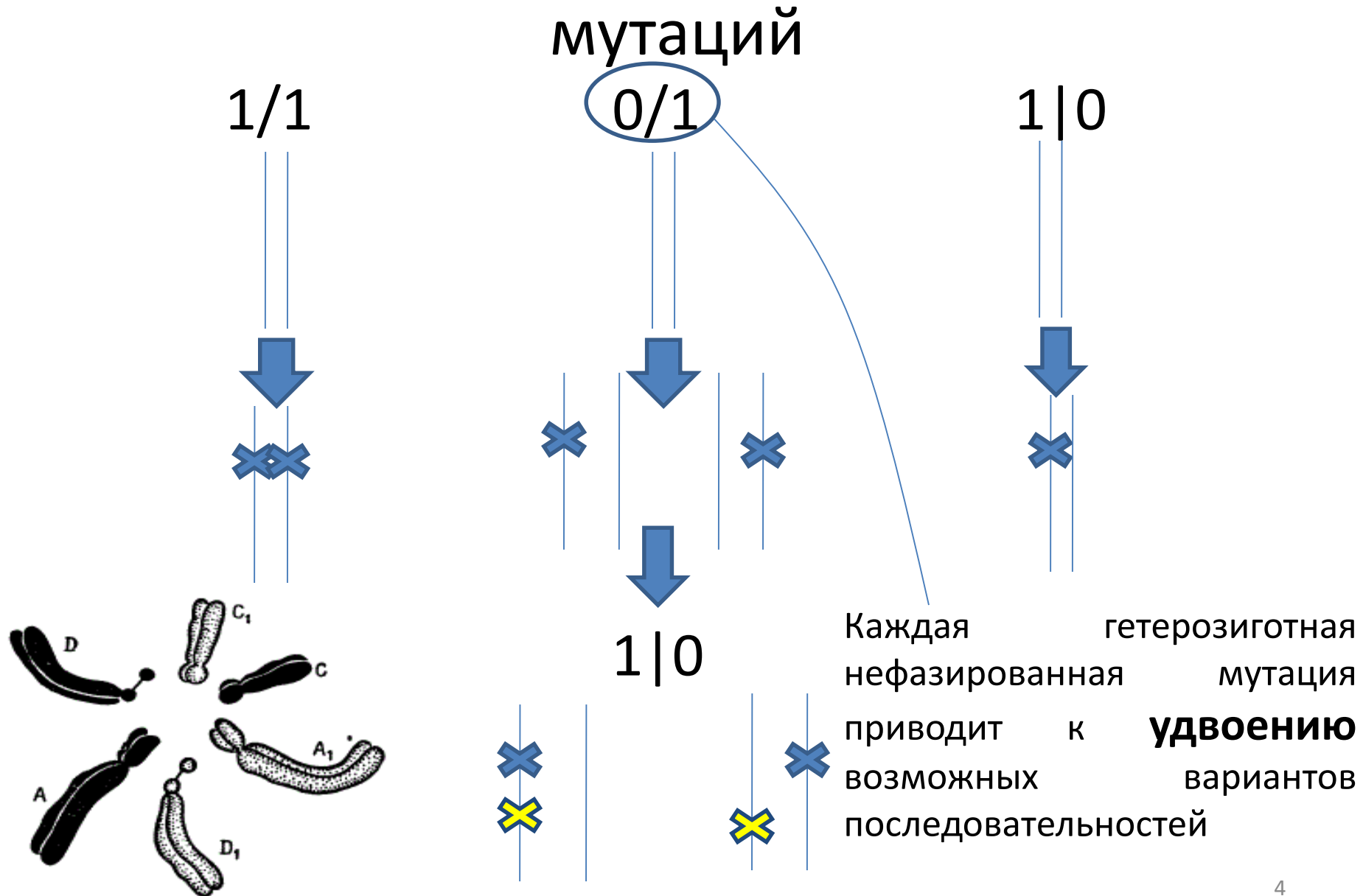


Задачи

- Изучение форматов файлов NGS, используемых в эксперименте (FASTA, VCF, GTF)
- Применение мутаций на референсную последовательность
- Восстановление аминокислотной последовательности с учетом мутаций



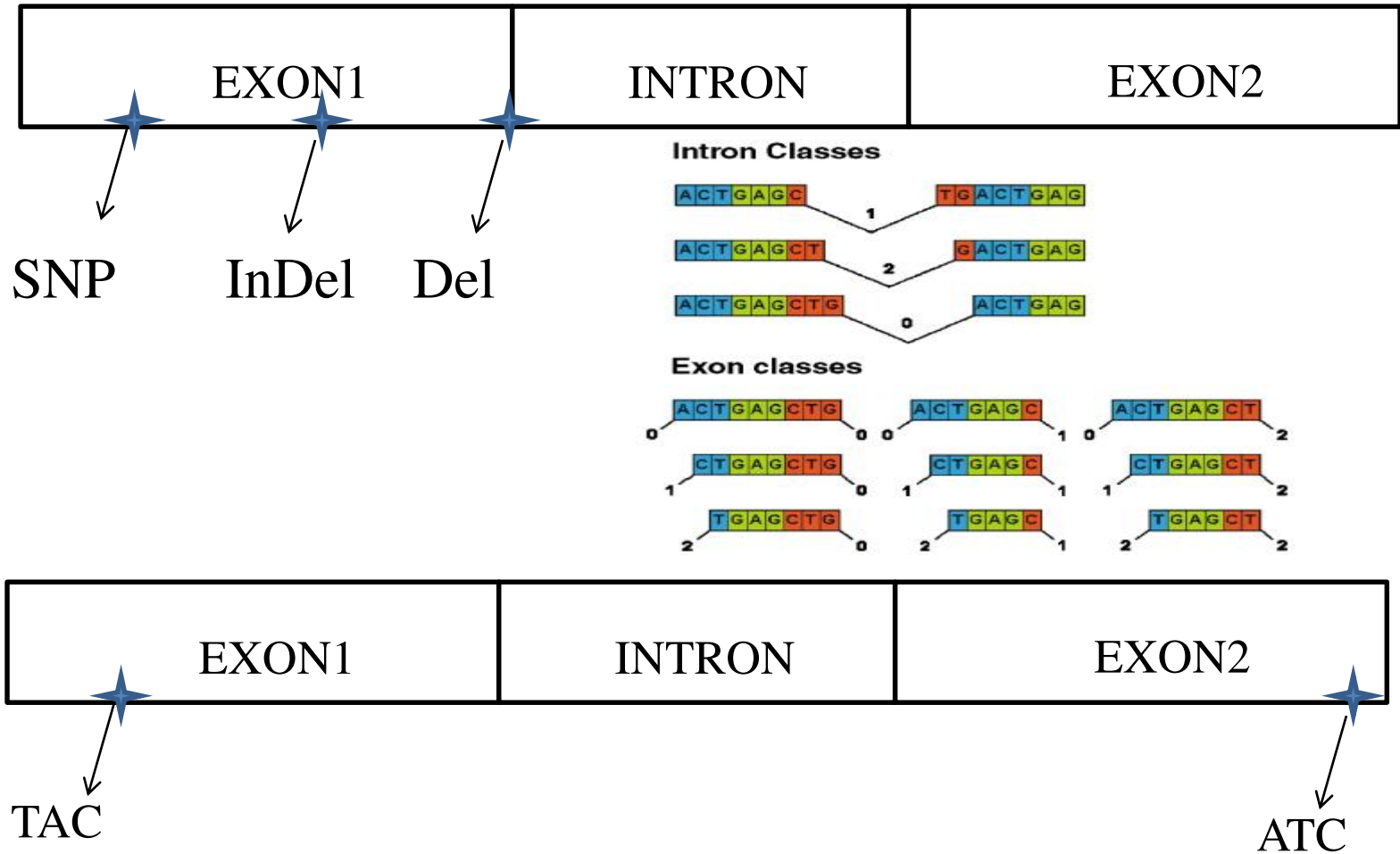
Различные варианты локализации мутаций



Входные данные

- Референсная последовательность дрозофилы (Fasta)
- Аннотации генов (GTF)
- Список мутаций (VCF) – 30 мутаций, из них:
 - 15 гетерозиготных, нефазированных, из них:
 - 10 в экзонах, из них:
 - 6 в белок-кодирующей части
- Таким образом, гипотетическое пространство вариантов равно $2^6 = 64$

Изменение интрон/экзонной картины

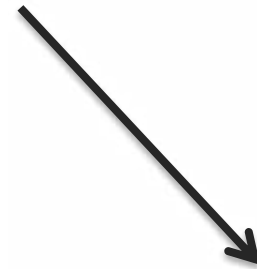


Подводные камни

- Влияние мутаций друг на друга
- Изменение положений экзонов и стартового кодона при мутациях
- Направление считывания



Языки, используемые в работе:



Python

Java

R

- Быстрый прототип
- Библиотека biopython

- Скорость работы
- Библиотека htsjdk

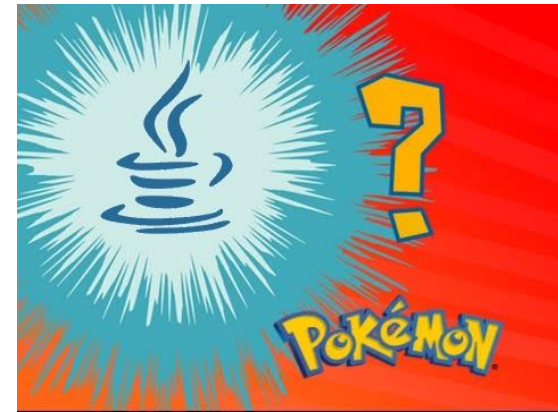
- Построение графиков
- Библиотека bioconductor

Библиотека ProRecon на языке Java

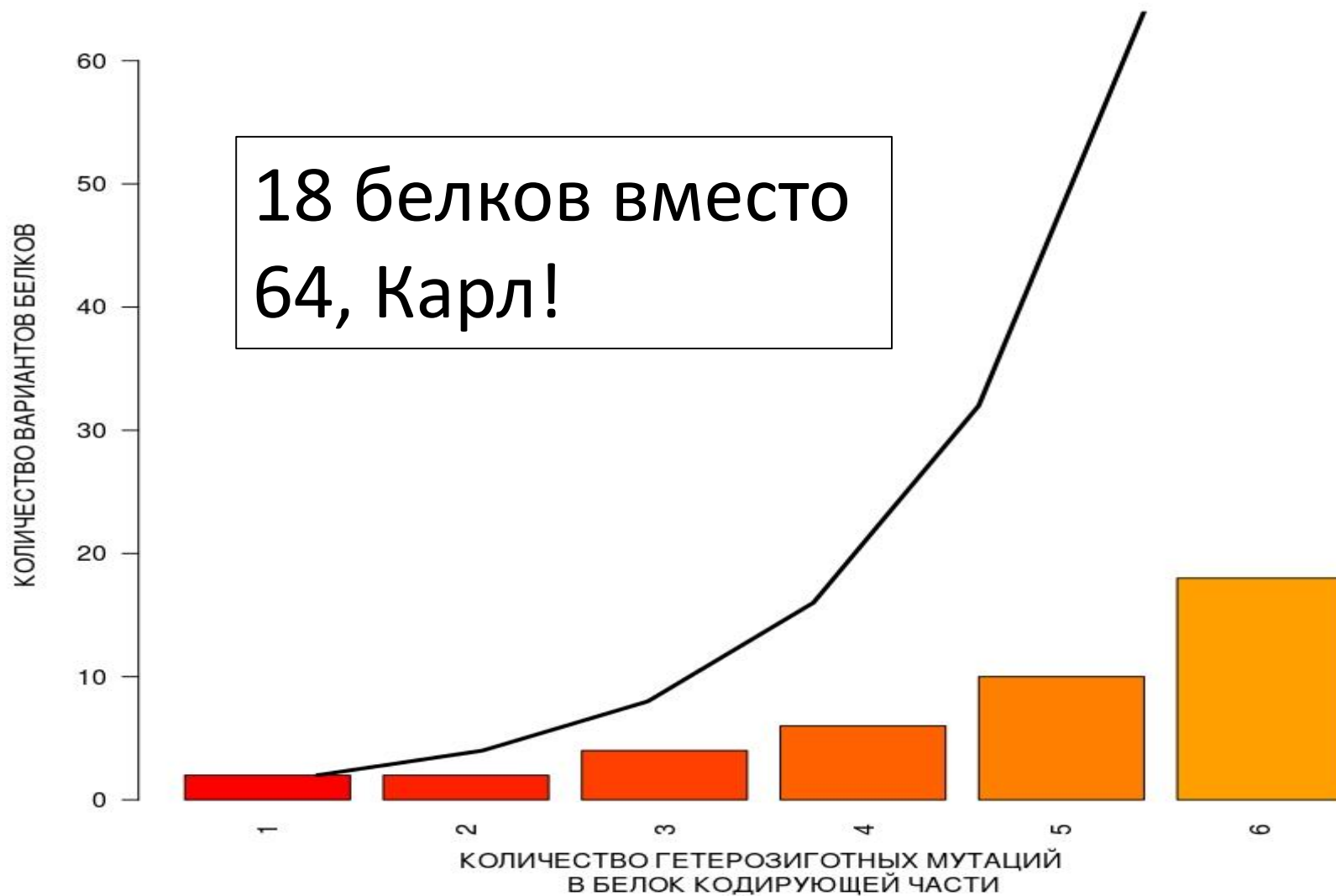
- Protein Reconstruction → ProRecon
- Полностью функциональна
- Выложена вместе с jar на github
- Покрыта тестами
- Качество существенно повышено за счёт

сравнения результатов с полученными в python

https://github.com/Petr1990/protein_reconstruction



Количество гипотетических и возможных белков



Произшедшие изменения белка

Делеция

```
RNLTPDTESKERALKKALK
RNLTPDT-----ALKKALK
RNLTPDT-----ALKKALK
RNLTPDTESKERALKKALK
RNLTPDTESKERALKKALK
```

Мутация перед
СТОП КОДОНОМ

```
VRTHFNTRC-----
VRTHFNTRLLARRSSHTLY
VRTHFNTRC-----
VRTHFNTRLLARRSSHTLY
VRTHFNTRC-----
```

Сдвиг рамки
считывания

```
PDFMPRNSD-----FSLNQ
PDFMCLAIRTSV-----FSLNQ
PDFMCLAIRTSV-----FSLNQ
PDFMPRNSD-----FSLNQ
PDFMPRNSD-----FSLNQ
```

Мутация провоцирующая
СТОП КОДОН

```
GEHKFHP ECF C C T A C G S F I G I
GEHKFHP ECF -----
GEHKFHP ECF C C T A C G S F I G I
GEHKFHP ECF C C T A C G S F I G I
GEHKFHP ECF C C T A C G S F I G I
```

Планы на будущее

1) Учёт сложных случаев

- Совпадение координат мутаций
- Учёт структурных вариаций (инверсия, БНД)
- Обработка мутаций затрагивающих сайты сплайсинга

2) Доработка качества библиотеки для полноценного релиза на гитхаб и хранилище Java библиотек

3) Интеграция в геномный браузер NGB

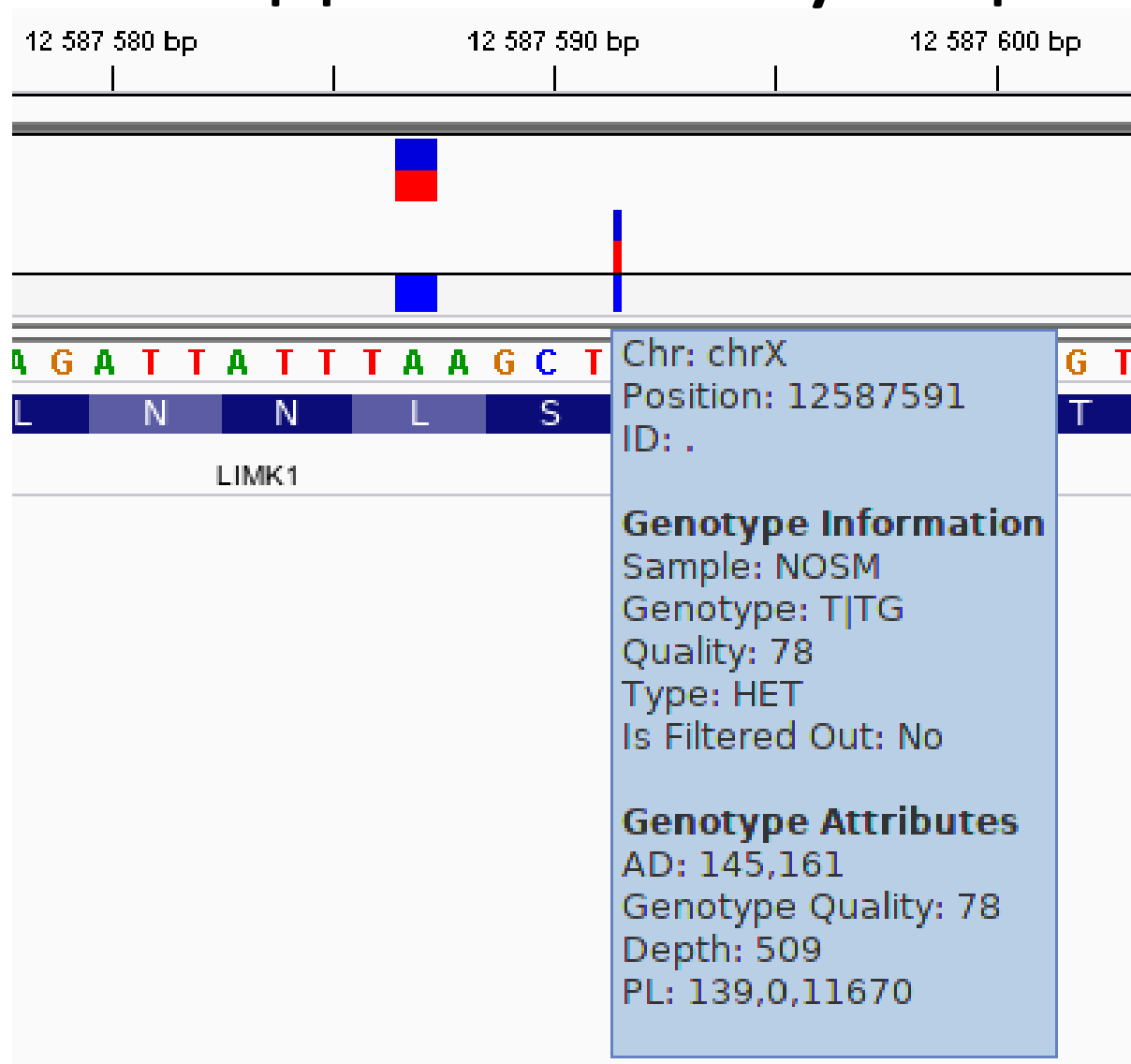
4)

5) PROFIT!!!

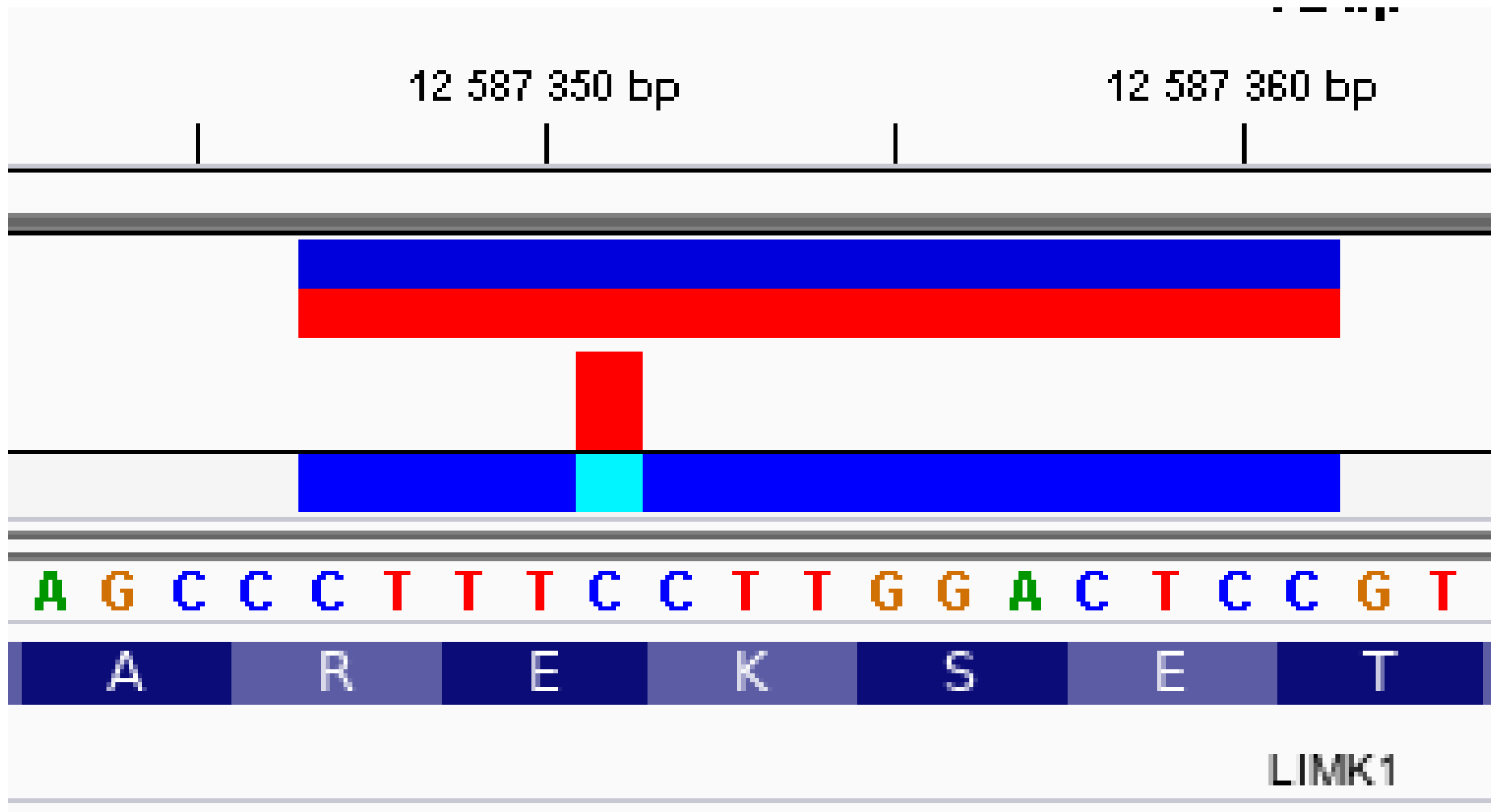
Спасибо за внимание!



Взаимодействие мутаций



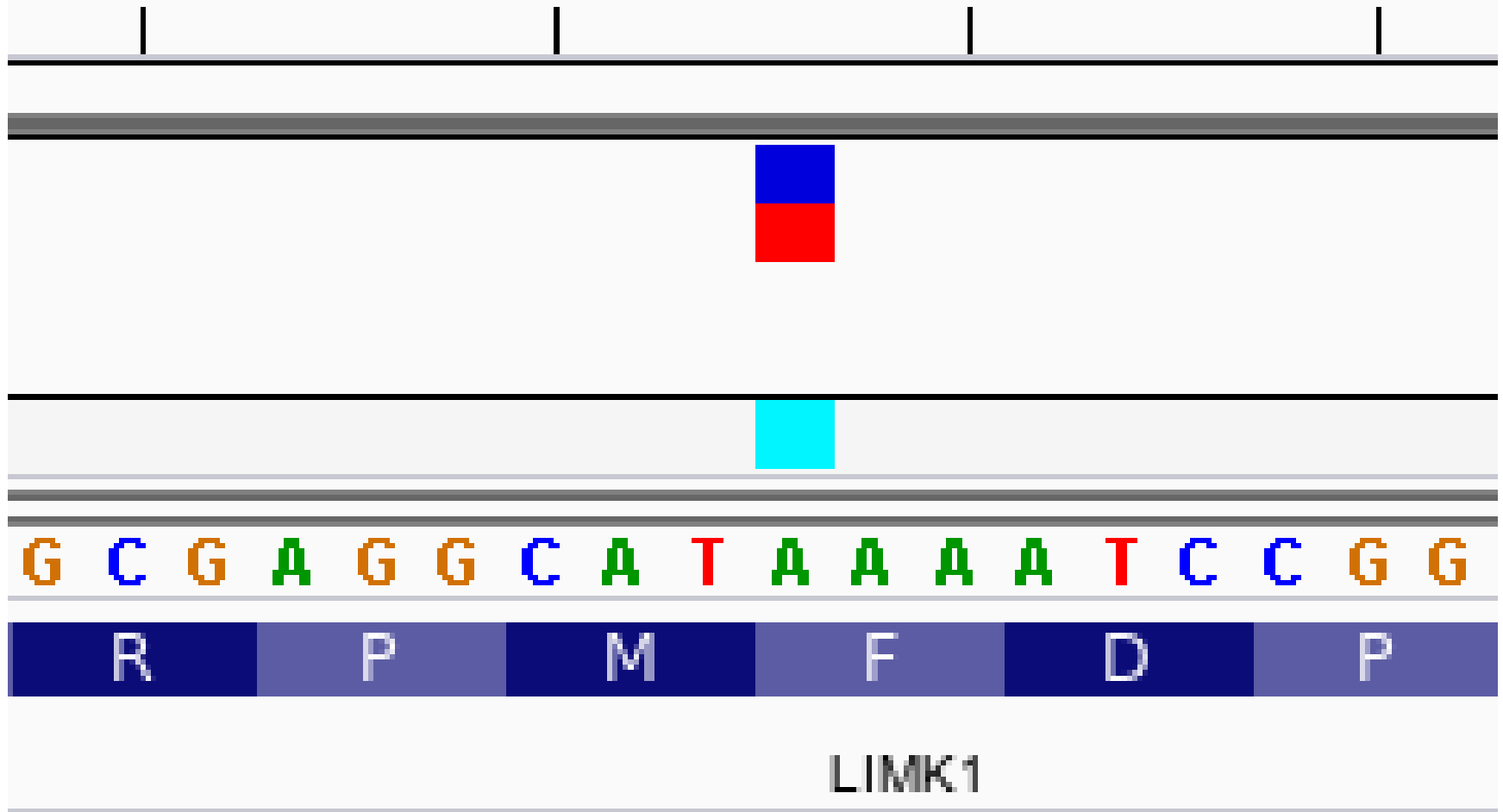
Наложение мутаций



Пример гетерозиготной фрейм-шифт делеции

2 587 860 bp

12 587 870 bp



Попытка выравнивания фрейм-шифта

```
sel=0 587 732
2
6 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDL-----CL-
257 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDL-----CL-
262 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
513 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
774 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
1025 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
1286 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
1537 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
1798 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
2049 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
2310 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
2561 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
2822 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
3073 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
3334 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
3585 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
3846 RORYTVVGNPYWMAPEMMKGLKYDEKVDVFSFGIMLCEIIGRVEADPDFMPRNSDFSLNQEFREKFCACPEPFVKVAFVCCDLNPDMRPCFETLHWLORLADDLAADRVPPE
```